

Available online at www.sciencedirect.com

Applied Mathematics Letters 20 (2007) 312–315

**Applied
Mathematics
Letters**www.elsevier.com/locate/aml

An improved class of estimators of a finite population quantile in sample surveys[☆]

A. Arcos^{*}, M. Rueda, J.F. Muñoz*Department of Statistics and Operational Research, University of Granada, Spain*

Received 28 July 2005; received in revised form 19 April 2006; accepted 24 April 2006

Abstract

This work proposes a general class of estimators for a finite population quantile using auxiliary information. This information is provided by the population means of auxiliary variables. The optimum estimator in this class is derived. This result is supported with a numerical example.

© 2006 Elsevier Ltd. All rights reserved.

Keywords: Auxiliary information; Difference type estimator; Finite population quantile; Horvitz–Thompson estimator; Sample surveys

1. Introduction

The utility of estimating population quantiles to describe the distribution of a population characteristic of interest is well known.

Most research into quantile estimation [1,5] deals exclusively with the interest variable and does not make explicit use of auxiliary variables in the construction of estimators. In practice, it is normal to have auxiliary information on one or more population characteristics related to the principal variable. In many situations, the information is associated with the interest variable but studied on a previous occasion when a census was carried out. Kuk and Mak [2], Rao et al. [3], Rueda et al. [4] and Singh et al. [7] have considered the problem of median estimation when the population median of an auxiliary variable is known.

In this work we suggest a class of estimators for a finite population quantile of the interest variable when the population means of the auxiliary variables are known.

2. A proposed class of estimators for the quantiles

Assume that a sample s of size n is drawn from a finite population U of size N by a specific sampling design d with first and second order positive inclusion probabilities, π_i and π_{ij} . Let y be the interest variable (which is the object of study) and x the auxiliary variable. For each $t \in \mathbb{R}$, the finite population distribution function of y is defined

[☆] Research partially supported by CICE (Junta de Andalucía) contract no. SEJ565.

^{*} Corresponding author. Tel.: +34 58243267; fax: +34 58249046.

E-mail address: arcos@ugr.es (A. Arcos).

by $F_y(t) = N^{-1} \sum_{i \in U} \Delta(t - y_i)$, with $\Delta(a) = 1$ if $a \geq 0$ and $\Delta(a) = 0$ otherwise. The finite population β -quantile, $Q_y(\beta)$ ($0 < \beta < 1$), is defined by $Q_y(\beta) = F_y^{-1}(\beta) = \inf\{y_i \in U : F_y(y_i) \geq \beta\}$.

First, assume that there are x_1, \dots, x_k auxiliary variables, and that the finite population β -quantiles $Q_{xi}(\beta)$, $i = 1, \dots, k$ are unknown, but that the finite population means $\bar{X}_1, \dots, \bar{X}_k$ of the auxiliary variables are known. The usual unbiased estimators (the Horvitz–Thompson estimators) of \bar{X}_i and $Q_y(\beta)$ are $\bar{x}_{HTi} = N^{-1} \sum_{j \in s} x_{ij}/\pi_j$ and $\hat{Q}_y(\beta) = \hat{F}_{HTy}^{-1}(\beta) = \inf\{y_i \in s : \hat{F}_{HTy}(y_i) \geq \beta\}$, respectively, with $\hat{F}_{HTy}(t) = N^{-1} \sum_{j \in s} \Delta(t - y_j)/\pi_j$.

Following the approach adopted by Srivastava and Jhajj [8], we propose a family of estimators of $Q_y(\beta)$ as

$$t_w = H(\hat{Q}_y(\beta), u_1, \dots, u_k), \quad (1)$$

where $u_i = \bar{x}_{HTi}/\bar{X}_i$, and $H(\cdot)$ is a function that is continuous and bounded in a closed subset $P \subset \mathbb{R}^{k+1}$ containing the point $(Q_y(\beta), \bar{X}_1, \dots, \bar{X}_k)$ such that:

- $H(Q_y(\beta), 1, \dots, 1) = Q_y(\beta)$ and the first order partial derivative of H with respect to $\hat{Q}_y(\beta)$, $H_0(Q_y(\beta), 1, \dots, 1) = 1$.
- The first and second order partial derivatives of $H(\hat{Q}_y(\beta), u_1, \dots, u_k)$ exist and are continuous and bounded in P .

Note that $\hat{Q}_y(\beta)$ and $\hat{Q}_{ri}(\beta) = \hat{Q}_y(\beta)\bar{X}_i/\bar{x}_{HTi}$, $\forall i$, are included in (1).

Expanding this function H about the point $(Q_y(\beta), 1, \dots, 1)$ in a second order Taylor series we obtain that the bias of t_w is of order $O(n^{-1})$ and the mean squared error of t_w , to the first degree of approximation, is given by

$$\text{MSE}(t_w) = V(\hat{Q}_y(\beta)) + H'CH + 2H'C_0,$$

where $H' = (H_1, \dots, H_k)$, H_i is the first order partial derivative of H with respect to u_i at the point $(Q_y(\beta), 1, \dots, 1)$, $C = (c_{ij})_{k \times k}$, $C_0 = (c_{01}, \dots, c_{0k})'$, $c_{0i} = \text{cov}(\hat{Q}_y(\beta), \bar{x}_{HTi})Q_y(\beta)/\bar{X}_i$ and $c_{ij} = \text{cov}(\bar{x}_{HTi}, \bar{x}_{HTj})Q_y(\beta)^2/(\bar{X}_i\bar{X}_j)$, $i = 1, \dots, k$.

Proposition 1. Up to terms of order $O(n^{-1})$,

$$\text{MSE}(t_w) \geq V(\hat{Q}_y(\beta)) - C_0'C^{-1}C_0. \quad (2)$$

Proof. The optimum value of H can be obtained by differentiating the above expression for the $\text{MSE}(t_w)$ and equating to zero. We obtain $H = -C^{-1}C_0$. \square

Exponential and difference type estimators $\hat{Q}_{y_{ex1}}(\beta) = \hat{Q}_y(\beta) \prod_{i=1}^k (\bar{X}_i/\bar{x}_{HTi})^{\alpha_i}$ and $\hat{Q}_{y_{dif1}}(\beta) = \hat{Q}_y(\beta) + \sum_{i=1}^k c_i(\bar{x}_{HTi} - \bar{X}_i)$, are also included in (1). The optimum values of α_i and c_i ($i = 1, \dots, k$) can be obtained by minimizing $\text{MSE}(t_w)$ with a respective function H . The resulting estimators with the optimum constants have the minimum MSE, which is given by (2). Thus, asymptotically, $\hat{Q}_{y_{ex1}}(\beta)$ and $\hat{Q}_{y_{dif1}}(\beta)$ are optimum in this class.

3. The proposed estimators under simple random sampling

Under SRSWOR design, \bar{x}_{HTi} reduces to the sample mean, $\bar{x}_i = n^{-1} \sum_{j \in s} x_{ij}$, $\hat{F}_{HTy}(t)$ reduces to the ordinary sample empirical distribution function, and $\hat{Q}_y(\beta) = \hat{F}_y^{-1}(\beta)$ reduces to the sample β -quantile of y . In this particular case, variances and covariances are easily obtained (see the Appendix):

$$\begin{aligned} v(\bar{x}_i) &= \frac{1-f}{n} S_{xi}^2, & v(\hat{Q}_y(\beta)) &= \frac{1-f}{n} \beta(1-\beta) \left(\frac{1}{f_y(Q_y(\beta))} \right)^2, \\ \text{cov}(\bar{x}_i, \bar{x}_j) &= \frac{1-f}{n} S_{xixj}, & \text{cov}(\bar{x}_i, \hat{Q}_y(\beta)) &= -\frac{1-f}{n} \frac{\rho_{zxi} S_{xi} \sqrt{\beta(1-\beta)}}{f_y(Q_y(\beta))}, \end{aligned}$$

where $z_j = \Delta(Q_y(\beta) - y_j)$, ρ_{zxi} denotes the coefficient of correlation between z and x_i . Variances and covariances can be estimated from the sample and $f_y(Q_y(\beta))$ can be approximated following Silverman [6].

Asymptotic properties of estimators in the proposed class are derived assuming that the finite population embeds in a sequence of populations $\{U_v\}$, where n_v and N_v increase such that $\frac{n_v}{N_v} \rightarrow f$ when $n_v \rightarrow \infty$. It is also assumed

that when $N_y \rightarrow \infty$, the multivariate distribution can be approximated by a continuous distribution with marginal densities f_y and f_{x_i} for y and x_i , ($i = 1, \dots, k$) respectively, and $f_{x_i}(Q_{x_i}(\beta))$ are positive.

Under simple random sampling, the estimators in the proposed class t_w are asymptotically unbiased and normal. First, the asymptotic unbiasedness of the proposed class of estimators is easily derived from its linear expression and since the estimators $\hat{Q}_y(\beta)$ and \bar{x}_i are, respectively, asymptotically unbiased and unbiased for their respective parameters. Second, as \bar{x}_i and $\hat{Q}_y(\beta)$ are asymptotically normal [1], then so are the estimators in t_w .

4. A more general class of estimators

We consider the situation in which the population means and the population β -quantiles associated with the auxiliary variables x_i , \bar{X}_i and $Q_{x_i}(\beta)$ are known. A family of estimators of $Q_y(\beta)$ that is broader than t_w is defined by

$$T_w = H(\hat{Q}_y(\beta), u_1, \dots, u_k, v_1, \dots, v_k), \quad (3)$$

where $u_i = \bar{x}_{HTi}/\bar{X}_i$, $v_i = \hat{Q}_{x_i}(\beta)/Q_{x_i}(\beta)$ and $H(\hat{Q}_y(\beta), u_1, \dots, u_k, v_1, \dots, v_k)$ is a function that satisfies similar regularity conditions to those given in Section 3. Eq. (3) provides a class of estimators, of which t_w is a particular case. This class also includes:

- the ratio estimator proposed by Kuk and Mak [2], $\hat{Q}_r(\beta) = \hat{Q}_y(\beta) \frac{Q_{x_i}(\beta)}{\hat{Q}_{x_i}(\beta)}$,
- the difference estimator proposed by Rao et al. [3],

$$\hat{Q}_{dr}(\beta) = \hat{Q}_y(\beta) + \hat{R} \left(Q_x(\beta) - \hat{Q}_x(\beta) \right), \quad \text{with } \hat{R} = \frac{\sum_{i \in s} y_i / \pi_i}{\sum_{i \in s} x_i / \pi_i},$$

- the regression estimator proposed by Rueda et al. [4],

$$\hat{Q}_{reg}(\beta) = \hat{Q}_y(\beta) + \hat{b} \left(Q_x(\beta) - \hat{Q}_x(\beta) \right), \quad \text{with } \hat{b} = \frac{\sum_{i \in s} x_i y_i / \pi_i}{\sum_{i \in s} x_i^2 / \pi_i}.$$

We can write (3) as $T_w = H(\hat{Q}_y(\beta), w_1, \dots, w_{2k})$ with $w_i = u_i$ for $i = 1, \dots, k$ and $w_i = v_i$ for $i = k + 1, \dots, 2k$. The bias and the MSE of the estimators in the class T_w can be derived using a method similar to the one followed in the previous section.

Exponential and difference estimators:

$$\begin{aligned} \hat{Q}_{y_{ex2}}(\beta) &= \hat{Q}_y(\beta) \prod_{i=1}^k \left(\frac{\bar{X}_i}{\bar{x}_{HTi}} \right)^{\gamma_i} \prod_{i=1}^k \left(\frac{Q_{x_i}(\beta)}{\hat{Q}_{x_i}(\beta)} \right)^{\delta_i}, \\ \hat{Q}_{y_{dif2}}(\beta) &= \hat{Q}_y(\beta) + \sum_{i=1}^k a_i (\hat{Q}_{x_i}(\beta) - Q_{x_i}(\beta)) + \sum_{i=1}^k d_i (\bar{x}_{HTi} - \bar{X}_i), \end{aligned}$$

are also asymptotically optimum in T_w . Thus, asymptotically, $\hat{Q}_{y_{ex2}}(\beta)$ and $\hat{Q}_{y_{dif2}}(\beta)$ are more efficient than all the cited estimators.

5. Numerical comparisons

The following examples reflect the potential gains from the use of the proposed estimators instead of the customary estimators. Three well known populations are considered (see [4] for details): SUGAR CANE, MU281 and FAM1500.

Let consider us SRSWOR and the median estimation, $Q_y(0.5) = M_y$. Assuming that the mean or the median and the median associated with an auxiliary variable, \bar{X} and M_x , are known, the following five estimators are compared: the simple estimator, \hat{M}_y , the usual difference (\hat{M}_d) and ratio (\hat{M}_r) estimators and the proposed estimators $\hat{M}_{D1} = \hat{M}_y + \hat{c}_{opt}(\bar{X} - \bar{x})$, and $\hat{M}_{D2} = \hat{M}_y + \hat{a}_{opt}(M_x - \hat{M}_x) + \hat{d}_{opt}(\bar{X} - \bar{x})$. The root of the ratio, R , of the variance to the variance of the simple estimator is computed as a measure of efficiency, $R(\hat{\theta}) = \sqrt{V(\hat{\theta})/V(\hat{M}_y)}$.

Table 1
Relative efficiencies (R) of the various estimators when the median is estimated

	\hat{M}_y	\hat{M}_d	\hat{M}_r	\hat{M}_{D1}	\hat{M}_{D2}
SUGAR CANE	1.0000	0.7882	1.1848	0.8058	0.7525
MU281	1.0000	0.9125	0.9317	0.8781	0.8766
FAM1500	1.0000	0.7477	0.8023	0.7074	0.6912

From Table 1 we deduce that the proposed estimators are more efficient than the simple estimator which does not use information of the population means. Evidently, \hat{M}_{D2} is the most accurate since this estimator uses population information (mean and median) from the auxiliary variable.

Appendix

In this appendix, $\text{cov}(\bar{x}_i, \hat{Q}_y(\beta))$ is derived. The Taylor series expansion yields

$$\hat{Q}_y(\beta) - Q_y(\beta) \simeq \frac{1}{f_y(Q_y(\beta))} (\hat{F}_y(\hat{Q}_y(\beta)) - \hat{F}_y(Q_y(\beta))) + O(n^{-\frac{1}{2}}) \simeq \frac{(\beta - \bar{z}_y)}{f_y(Q_y(\beta))}$$

where $\bar{z}_y = \hat{F}_y(Q_y(\beta)) = n^{-1} \sum_{j \in s} z_j$ and $f_y(\cdot)$ is the derivative of $\hat{F}_y(\cdot)$, the limiting value of $F_y(\cdot)$ as $N \rightarrow \infty$. Then

$$\begin{aligned} E(\bar{x}_i - \bar{X}_i)(\hat{Q}_y(\beta) - Q_y(\beta)) &\simeq \frac{1}{f_y(Q_y(\beta))} E(\bar{x}_i - \bar{X}_i)(\beta - \bar{z}_y) \\ &= -\frac{\text{cov}(\bar{x}_i, \bar{z}_y)}{f_y(Q_y(\beta))} = -\frac{1-f}{n} \frac{\rho_{zxi} S_{xi} S_z}{f_y(Q_y(\beta))} = -\frac{1-f}{n} \frac{\rho_{zxi} S_{xi} \sqrt{\beta(1-\beta)}}{f_y(Q_y(\beta))}. \end{aligned}$$

References

- [1] S.T. Gross, Median estimation in sample survey, Proc. Surv. Res. Meth. Sect. Amer. Statist. Ass. (1980) 181–184.
- [2] A. Kuk, T.K. Mak, Median estimation in the presence of auxiliary information, J. Roy. Statist. Soc. Ser. B 1 (1989) 261–269.
- [3] J.N.K. Rao, J.G. Kovar, H.J. Mantel, On estimating distribution functions and quantiles from survey data using auxiliary information, Biometrika 77 (1990) 365–375.
- [4] M. Rueda, A. Arcos, M.D. Martínez, Difference estimators of quantiles in finite populations, Test 12 (2003) 481–496.
- [5] J. Sedransk, P.J. Smith, Inference for finite population quantiles, in: P.R. Krishnaiah, C.R. Rao (Eds.), Handbook of Statistics, vol. 6, North-Holland, 1988, pp. 267–289 (Chapter 11).
- [6] B.W. Silverman, Density Estimation for Statistics and Data Analysis, Chapman and Hall, London, 1986.
- [7] S. Singh, A.H. Joarder, D.S. Tracy, Median estimation using double sampling, Aust. N. Z. J. Stat. 43 (2001) 33–46.
- [8] S.K. Srivastava, H.S. Jhajj, A class of estimators of the population mean in survey sampling using auxiliary information, Biometrika 68 (1981) 341–343.